



**POLITECNICO**  
**MILANO 1863**

SCUOLA DI INGEGNERIA INDUSTRIALE  
E DELL'INFORMAZIONE

# SMBUD Project - Spark

Author(s): **Gabriele Ginestroni**

**Giacomo Gumiero**

**Lorenzo Iovine**

**Nicola Landini**

**Francesco Leone**

Group Number: **10**

Academic Year: 2022-2023



# Contents

Contents	i
<b>1 Introduction</b>	<b>1</b>
<b>2 Data Structure</b>	<b>3</b>
2.1 Article structure . . . . .	3
2.2 Author structure . . . . .	5
<b>3 Data Import</b>	<b>7</b>
3.1 Setup . . . . .	7
3.2 Article Custom Schema . . . . .	7
3.3 Venue Collection Structure . . . . .	8
3.4 Author Collection Structure . . . . .	9



# 1 | Introduction

In this chapter will be presented the problem specification and the hypothesis under which the database is implemented.



## 2 | Data Structure

### 2.1. Article structure

In the following JSON file we can see how the *Article* is structured in our dataset.

**Note:** in some fields we decided not to insert all the content only for ease of read reasons.

```

1 {
2   "_id": "53e99f86b7602d9702859fdf",
3   "title": "Locality Sensitive Outlier Detection: A ranking driven
      approach",
4   "authors": [
5     {
6       "idAuth": "542a4c9fdabfae61d496694e",
7       "org": "Computer Science and Engineering Department, The Ohio State
      University, USA"
8     },
9     {
10      "idAuth": "53f48bc5dabfaea7cd1cce1d",
11      "org": "Computer Science and Engineering Department, The Ohio State
      University, USA"
12     },
13     {
14      "idAuth": "53f44b6fdabfaec09f1dd00d",
15      "org": "Computer Science and Engineering Department, The Ohio State
      University, USA"
16     }
17   ],
18   "n_citation": 60,
19   "abstract": "Outlier detection is fundamental to a variety of database
      and analytic tasks. Recently, distance-based outlier detection has
      emerged as a viable and scalable alternative to traditional
      statistical and geometric approaches. In this article we explore the
      role of ranking for the efficient discovery of distance-based
      outliers from large high dimensional data sets. Specifically, we
      develop a light-weight ranking scheme that is powered by locality

```

sensitive hashing, which reorders the database points according to their likelihood of being an outlier. We provide theoretical arguments to justify the rationale for the approach and subsequently conduct an extensive empirical study highlighting the effectiveness of our approach over extant solutions. We show that our ranking scheme improves the efficiency of the distance-based outlier discovery process by up to 5-fold. Furthermore, we find that using our approach the top outliers can often be isolated very quickly, typically by scanning less than 3% of the data set.",

```
20 "doi": "10.1109/ICDE.2011.5767852",
21 "keywords": [
22     "database point",
23     "ranking scheme",
24     "geometric approach",
25     ...
26 ],
27 "isbn": "978-1-4244-8958-9",
28 "page_start": "410",
29 "page_end": "421",
30 "year": 2011,
31 "fos": [
32     "Locality-sensitive hashing",
33     "Anomaly detection",
34     "Data mining",
35     ...
36 ],
37 "venue": {
38     "raw": "ICDE",
39     "type": 0
40 },
41 "references": [
42     "53e99fddb7602d97028b7e65"
43 ]
44 }
```



## 2.2. Author structure

In the following JSON file we can see how the *Author* is structured in our dataset.

**Note:** in some fields we decided not to insert all the content only for ease of read reasons.

```
1 {  
2   "_id": "53f45775dabfaee4dc8162e6",  
3   "name": "Guillermo Jorge-Botana",  
4   "articles": ["53e99f86b7602d970285a187"],  
5   "orcid": "0000-0001-5879-6783",  
6   "bio": "Qing-Long Han received the B.Sc. degree in mathematics from the  
       ...",  
7   "email": "Guillermo.Jorge-Botana@yahoo.com",  
8   "nationality": "de",  
9   "dob": "1945-06-19T00:00:00Z"  
10 },
```



# 3 | Data Import

## 3.1. Setup

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.types import StructType, StructField, IntegerType,
   StringType, ArrayType
3 from pyspark.sql.functions import col, to_timestamp, rand, lit,
   collect_list, array
4
5 spark = SparkSession.builder \
6     .master("local") \
7     .appName("SparkByExamples.com") \
8     .getOrCreate()
```

## 3.2. Article Custom Schema

```
1 schemaArticle = StructType([
2     StructField('_id', StringType(), True),
3     StructField('title', StringType(), True),
4     StructField('authors',
5         ArrayType(
6             StructType([
7                 StructField('idAuth', StringType(), True),
8                 StructField('org', StringType(), True)
9             ]), True)
10 ),
11 StructField('n_citation', IntegerType(), True),
12 StructField('abstract', StringType(), True),
13 StructField('doi', StringType(), True),
14 StructField('keywords', ArrayType(StringType()), True),
15 StructField('isbn', StringType(), True),
16 StructField('page_start', StringType(), True),
17 StructField('page_end', StringType(), True),
18 StructField('year', IntegerType(), True),
19 StructField('fos', ArrayType(StringType()), True),
```

```

20     StructField('references', ArrayType(StringType()), True),
21     StructField('venue',
22         StructType([
23             StructField('raw', StringType(), True),
24             StructField('type', IntegerType(), True),
25             StructField('issue', StringType(), True),
26             StructField('volume', StringType(), True),
27             StructField('publisher', StringType(), True)
28         ])
29     ),
30 ])
```

We decided to use import from schema to explicitly show data structure:

```

1 df_articles = spark.read.schema(schemaArticle).json("../
                                dblp_sample_filtered_spark.json",
                                multiLine=True)
```

The attributes `issue`, `volume` and `publisher` inside `venue` are moved back in the root structure and removed from the inner struct:

```

1 df_articles = df_articles.withColumn("issue", col("venue.issue")) \
2                             .withColumn("volume", col("venue.volume")) \
3                             .withColumn("publisher", col("venue.publisher")) \
4                             .withColumn("venue", col("venue").dropFields("
                                issue", "volume", "publisher"))
```

### 3.3. Venue Collection Structure

A new dataframe is created with attributes of venue and the `_id` of the article, then it is all grouped by venue attributes and, for each venue, a list of the articles ids is created. Finally we drop rows with null `raw` to delete inconsistent tuple.

```

1 df_venues = df_articles.select("venue.raw", "venue.type", "_id") \
2                             .groupBy("raw", "type") \
3                             .agg(collect_list("_id").alias("artIds")) \
4                             .dropna(subset=["raw"])
```

Now we can keep only the `raw` attribute of the venue:

```
1 df_articles = df_articles.withColumn("venue_raw", col("venue.raw")).drop(
    ("venue"))
```

In the end, we add a generated field inside venues collection: for each venue a random city is selected that should represent the place where the venue was held.

```
1 citiesList = ["New York", "London", "Paris", "Berlin", "Madrid", "Rome",
    "Dublin", "Copenhagen", "Vienna", "
    Amsterdam", "Brussels", "Lisbon", "
    Prague", "Athens", "Budapest", "
    Warsaw", "Zurich", "Luxembourg", "
    Oslo", "Stockholm", "Helsinki", "
    Moscow", "Istanbul", "Kiev", "Minsk"
    , "Belgrade", "Bucharest", "Sofia",
    "Tallinn", "Riga", "Vilnius", "
    Tbilisi", "Yerevan", "Baku", "Dubai"
    , "Abu Dhabi", "Doha", "Manama", "
    Muscat", "Riyadh", "Jeddah", "Mecca"
    , "Medina", "Kuala Lumpur", "
    Singapore", "Hong Kong", "Shanghai",
    "Beijing", "Tokyo", "Seoul", "
    Bangkok", "Manila"]
2 cities = array([lit(city) for city in citiesList])
3 df_venues = df_venues.withColumn("city", cities.getItem((rand() * len(
    citiesList)).cast("int")))
```

### 3.4. Author Collection Structure

In order to handle authors' collection, we simply import from json with autogenerated schema, applying a conversion from string to timestamp.

```
1 df_authors = spark.read.json("../dblp_sample_reverted_filtered_spark.
    json", multiLine=True)
2 df_authors = df_authors.withColumn("dateofbirth", to_timestamp(
    df_authors["dob"], "yyyy-MM-dd'T'HH:
    mm:ss'Z'")) \
3     .drop("dob") \
4     .withColumnRenamed("dateofbirth", "dob")
```